

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6520108号  
(P6520108)

(45) 発行日 令和1年5月29日(2019.5.29)

(24) 登録日 令和1年5月10日(2019.5.10)

(51) Int.Cl.

F I

G 1 O L 13/10 (2013.01)

G 1 O L 13/10 1 1 2 E

G 1 O L 13/10 1 1 4

請求項の数 10 (全 14 頁)

(21) 出願番号 特願2014-259485 (P2014-259485)  
 (22) 出願日 平成26年12月22日(2014.12.22)  
 (65) 公開番号 特開2016-118722 (P2016-118722A)  
 (43) 公開日 平成28年6月30日(2016.6.30)  
 審査請求日 平成29年12月19日(2017.12.19)

(73) 特許権者 000001443  
 カシオ計算機株式会社  
 東京都渋谷区本町1丁目6番2号  
 (74) 代理人 100074099  
 弁理士 大菅 義之  
 (72) 発明者 田中 飛雄太  
 東京都羽村市栄町3丁目2番1号 カシオ  
 計算機株式会社羽村技術センター内  
 審査官 上田 雄

最終頁に続く

(54) 【発明の名称】 音声合成装置、方法、およびプログラム

(57) 【特許請求の範囲】

【請求項1】

テキストデータに基づいて、音声素片を選択する選択部と、  
 前記選択された音声素片を接続することにより接続音声素片を生成する接続部と、  
音声データから得られるピッチ列を量子化し、前記量子化されたピッチ列を平滑化し、  
前記平滑化されたピッチ列から抑揚情報を抽出する抑揚情報抽出部と、  
 前記接続音声素片に含まれる要素の少なくとも一部を前記抑揚情報に基づいて変更して  
 合成音声出力する出力部と、  
 を備える音声合成装置。

【請求項2】

前記出力部は、前記接続音声素片に含まれる要素であるピッチ列を、前記抑揚情報抽出  
 部から抽出された抑揚情報に適應するように変更するピッチ適應部を含む、請求項1に記  
 載の音声合成装置。

【請求項3】

前記平滑化されたピッチ列は、加重移動平均演算することにより平滑化されている、請  
 求項1または2に記載の音声合成装置。

【請求項4】

前記ピッチ適應部はさらに、前記音声データに含まれる韻律情報としてのピッチ列と前  
 記接続音声素片に含まれるピッチ列との時間スケールを調整するとともに、前記韻律情報  
 としてのピッチ列と前記接続音声素片に含まれるピッチ列とのピッチ存在区間を調整する

10

20

、請求項 2 に記載の音声合成装置。

【請求項 5】

前記出力部は、前記接続音声素片に含まれる要素であるパワー列を、前記抑揚情報抽出部から抽出された抑揚情報に適應するように変更するパワー適應部を含む、請求項 1 に記載の音声合成装置。

【請求項 6】

前記抑揚情報抽出部は、前記音声データに含まれる韻律情報としてのパワー列を平滑化し、当該平滑化されたパワー列を前記抑揚情報として抽出し、

前記パワー適應部は、前記接続音声素片に含まれるパワー列を平滑化し、当該平滑化されたパワー列と前記抑揚情報としての平滑化されたパワー列との比の列を算出し、当該比の列に基づいて前記接続音声素片のパワー列を修正する、請求項 5 に記載の音声合成装置。

10

【請求項 7】

前記平滑化されたパワー列は、前記パワー列に含まれるパワーそれぞれを加重移動平均演算することにより取得する、請求項 6 に記載の音声合成装置。

【請求項 8】

前記パワー適應部はさらに、前記音声データに含まれる韻律情報としてのパワー列及び前記接続音声素片に含まれるパワー列それぞれの時間スケールを調整する、請求項 5 ないし 7 のいずれかに記載の音声合成装置。

【請求項 9】

20

音声合成装置に用いられる音声合成方法であって、前記音声合成装置が、  
テキストデータに基づいて、音声素片を選択し、  
前記選択された音声素片を接続することにより接続音声素片を生成し、  
音声データから得られるピッチ列を量子化し、前記量子化されたピッチ列を平滑化し、  
前記平滑化されたピッチ列から抑揚情報を抽出し、

前記接続音声素片に含まれる要素の少なくとも一部を前記抑揚情報に基づいて変更して合成音声を出力する、音声合成方法。

【請求項 10】

音声合成装置に用いられる音声合成装置として用いられるコンピュータに、  
テキストデータに基づいて、音声素片を選択するステップと、  
前記選択された音声素片を接続することにより接続音声素片を生成するステップと、  
音声データから得られるピッチ列を量子化し、前記量子化されたピッチ列を平滑化し、  
前記平滑化されたピッチ列から抑揚情報を抽出するステップと、

30

前記接続音声素片に含まれる要素の少なくとも一部を前記抑揚情報に基づいて変更して合成音声を出力するステップと、  
を実行させるプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音声コーパスからの音声素片の選択によって音声合成を行う技術に関する。

40

【背景技術】

【0002】

入力テキストデータから生成される合成目標に対して、電子化された大規模な言語・音声データである音声コーパスを参照することにより音声波形の素片（以下、「音声素片」と記載する）を選択し、当該音声素片を接続することにより合成音声を出力する音声合成技術が知られている（例えば非特許文献 1～3 に記載の技術）。

【0003】

このような音声合成技術において、音声コーパスから合成目標に最も適合する音声素片列を選択するための手法として従来、次のような技術が知られている（例えば非特許文献 3 に記載の技術）。まず、入力テキストデータから抽出される音素列ごとに、その音素列

50

と同じ音素列を有する音声素片のデータ（以下、「素片データ」と記載する）が、素片候補データとして音声コーパスから抽出される。次に、DP（Dynamic Programming：動的計画法）アルゴリズムによって、入力テキストデータ全体に渡ってコストが最小となる最良の素片候補データの組（最良の素片データ列）が決定される。コストとしては、入力テキストデータと音声コーパス内の各素片データ間の音素列および韻律の差異、素片候補データである隣接する素片データ間のスペクトラム包絡などの音響パラメータ（特徴量ベクトルデータ）の不連続性などが用いられる。

#### 【0004】

入力テキストデータに対応する音素列は、例えば入力テキストデータに対して形態素解析処理を実行することで得られる。

10

#### 【0005】

入力テキストデータに対応する韻律（以下これを「目標韻律」と記載する）は、音素ごとの声帯の基本周波数であるピッチの高さ、持続時間長、および強度（パワー）である。この目標韻律の指定方式としては、入力テキストデータから得られる言語情報をもとに、実際の音声データに基づく統計的なモデルを用いて生成する方法がある（例えば非特許文献4に記載の技術）。言語情報は、例えば入力テキストデータに対して形態素解析処理を実行することで得られる。あるいは、目標韻律の指定方式として、ユーザが数値でパラメータ入力する方法がある。

#### 【0006】

さらに、目標韻律の指定方式として、ユーザが自身でテキストを発声するなどして与える音声によって指定するという方法がある。この方式は、テキストからの推定や数値パラメータの調整と比較して直感的な操作が可能であり、感情や抑揚の付与など自由度の高い目標韻律指定が行えるという利点がある。

20

#### 【先行技術文献】

#### 【非特許文献】

#### 【0007】

【非特許文献1】“CHATR:自然音声波形接続型任意音声合成システム”、電子情報通信学会信学技法、SP96-7.

【非特許文献2】“大規模コーパスを用いた音声合成システムXIMERA”、電子情報通信学会論文誌D Vol.J89-D No.12 pp.2688-2698.

30

【非特許文献3】河井 恒、“知識ベース 3-4 コーパスベース音声合成”、[online]、ver.1/2011.1.7、電子情報通信学会、[平成26年12月5日検索]、インターネット<URL: [http://27.34.144.197/files/02/02gun\\_07hen\\_03.pdf#page=6](http://27.34.144.197/files/02/02gun_07hen_03.pdf#page=6)>

【非特許文献4】匂坂 芳典、“知識ベース 3-7 韻律の生成”、[online]、ver.1/2011.1.7、電子情報通信学会、[平成26年12月5日検索]、インターネット<URL: [http://27.34.144.197/files/02/02gun\\_07hen\\_03.pdf#page=13](http://27.34.144.197/files/02/02gun_07hen_03.pdf#page=13)>

#### 【発明の概要】

#### 【発明が解決しようとする課題】

#### 【0008】

しかし、ユーザの音声入力による目標韻律指定方式には以下の様な課題がある。まず、指定できる目標韻律の自由度が増すため、それに対応できるだけの音声素片が必要となり、十分な量を揃えようとすると音声コーパスのデータベースが巨大になってしまう。また、ユーザが入力した音声の目標韻律と音声データベース内の音声素片が持っている韻律とでは、例えば声の高さなど個人に依存する相違点があり、適切な音声素片を選択することが困難になる。

40

#### 【0009】

上記の課題を解決するために、音声波形接続処理時に信号処理によって、音声素片の以下の要素を補正し、ユーザの入力した音声の目標韻律に適應させる方法が知られている。

1. 各音素の継続時間長。
2. ピッチ（音の高低）。

50

### ３．パワー（音の大小）。

しかし、信号処理によってユーザが入力した音声の目標韻律を音声データベースから選択した音声素片に単純に適應するだけでは、以下の様な問題が生じる。ユーザが入力した音声の目標韻律には細かなピッチやパワーの変動が含まれていることがあり、それら全てを音声素片に適應させると信号処理による音質の劣化が顕著になってしまう。また、ユーザの入力した音声の目標韻律と音声素片の韻律（特にピッチ）が大きく異なる場合、単純に適應してしまうと合成音声の音質が劣化してしまう。

#### 【００１０】

そこで、本発明は、波形接続方式の音声合成システムにおいて目標韻律を音声入力によって指定する際に、音声コーパスの規模拡大させる必要なく、高い自由度を維持しつつ合成音声の音質を向上させることを目的とする。

#### 【課題を解決するための手段】

#### 【００１１】

態様の一例では、入力されたテキストデータに基づいて、音声素片を選択する選択部と、選択された音声素片を接続することにより接続音声素片を生成する接続部と、入力された音声データに含まれる韻律情報から抑揚情報を抽出する抑揚情報抽出部と、接続音声素片に含まれる要素の少なくとも一部を抑揚情報に基づいて変更して合成音声を出力する出力部と、を備える。

#### 【発明の効果】

#### 【００１２】

本発明によれば、波形接続方式の音声合成システムにおいて目標韻律を音声入力によって指定する際に、音声コーパスの規模拡大させる必要なく、高い自由度を維持しつつ合成音声の音質を向上させることが可能となる。

#### 【図面の簡単な説明】

#### 【００１３】

【図１】音声合成装置の実施形態のブロック図である。

【図２】音声ＤＢのデータ構成例を示す図である。

【図３】音声合成装置の実施形態のハードウェア構成例を示す図である。

【図４】音声合成処理の例を示すフローチャートである。

【図５】ピッチ適應処理の説明図である。

【図６】パワー適應処理の説明図である。

【図７】ピッチ適應処理の詳細例を示すフローチャートである。

【図８】パワー適應処理の詳細例を示すフローチャートである。

#### 【発明を実施するための形態】

#### 【００１４】

以下、本発明を実施するための形態について図面を参照しながら詳細に説明する。図１は、音声合成装置の実施形態１００のブロック図である。音声合成装置１００は、音声合成部１０１、音声データベース（以下「音声ＤＢ」と記載）１０２、入力部１０３、および出力部１０４を備える。さらに、音声合成部１０１は、テキスト解析モジュール１０５、韻律解析モジュール１０６、素片選択モジュール１０７、波形接続モジュール１０８、ピッチ適應モジュール１０９、パワー適應モジュール１１０、およびシステム制御部１１１を備える。また、入力部１０３は、音声入力装置１１２およびテキスト入力装置１１３を備える。出力部１０４は、音声出力装置１１４を備える。素片選択モジュール１０７および波形接続モジュール１０８は音声素片選択・接続部に対応し、ピッチ適應モジュール１０９およびパワー適應モジュール１１０は抑揚情報抽出部および抑揚適應部に対応する。

#### 【００１５】

入力部１０３のテキスト入力装置１１３は、入力テキストデータを入力する。また、入力部１０３の音声入力装置１１２は、入力音声データを入力する。

#### 【００１６】

音声合成部 101 は、テキスト入力装置 113 から入力される入力テキストデータから生成される合成目標に対して、音声 DB 102 に記憶されている音声素片の集合である音声コーパスを参照することにより音声素片を選択し、その音声素片を接続することにより接続音声素片を生成する。

#### 【0017】

図 2 は、図 1 の音声 DB 102 に記憶される音声コーパスのデータ構成例を示す図である。音声コーパスとしては、例えば下記項目のデータが格納されている。

- ・ 予め録音された音声データ（図 2（a））。
- ・ 図 2（a）の音声データに対する音素ラベルの情報（図 2（b））。基本的にこのラベル付けされた図 2（a）の音声データの断片が音声素片となる。この音素ラベルの情報は、図 2（b）に示されるように、「開始位置」、「継続時間長」、および「音素種類」の各情報を有する。
- ・ 一定時間 T（ms：ミリ秒）からなるセグメントごとに、図 2（a）の音声データから解析されたピッチ、パワー、フォルマントなどの音響情報（図 2（c））。セグメント長 T は、例えば「10」ms である。

#### 【0018】

図 1 の説明に戻り、音声合成部 101 内のテキスト解析モジュール 105 は、テキスト入力装置 113 が入力した入力テキストデータに対して例えば形態素解析処理を実行することにより、入力テキストデータに対応するアクセント付きの音素列を抽出する。

#### 【0019】

音声合成部 101 内の韻律解析モジュール 106 は、音声入力装置 112 が入力した入力音声データを解析して、目標韻律を抽出する。

#### 【0020】

音声合成部 101 内の素片選択モジュール（音声素片選択・接続部）107 は、入力テキストデータから生成された音素列と入力音声データから生成された目標韻律とからなる合成目標に対して、音声データ内の音声コーパス（図 2）を参照することにより音声素片を選択する。

#### 【0021】

音声合成部 101 内の波形接続モジュール 108 は、素片選択モジュール 107 で選択された音声素片を接続することにより、接続音声素片を生成する。

#### 【0022】

音声合成部 101 内のピッチ適応モジュール 109 は、波形接続モジュール 108 が出力する接続音声素片に含まれるピッチ列を、入力部 103 の音声入力装置 112 から入力された入力音声データに含まれるピッチ列に適応させるように変更する。

#### 【0023】

音声合成部 101 内のパワー適応モジュール 110 は、波形接続モジュール 108 が出力する接続音声素片に含まれるパワー列を、入力部 103 の音声入力装置 112 から入力された入力音声データに含まれるパワー列に適応させるように変更する。

#### 【0024】

音声合成部 101 内のシステム制御部 111 は、音声合成部 101 内の 105 ~ 110 の各部分の動作の実行順序等を制御する。

#### 【0025】

図 3 は、図 1 の音声合成装置 100 をソフトウェア処理として実現できるコンピュータのハードウェア構成例を示す図である。図 3 に示されるコンピュータは、CPU 301、ROM（リードオンリーメモリ：読出し専用メモリ）302、RAM（ランダムアクセスメモリ）303、入力装置 304、出力装置 305、外部記憶装置 306、可搬記録媒体 410 が挿入される可搬記録媒体駆動装置 307、及び通信インターフェース 308 を有し、これらがバス 309 によって相互に接続された構成を有する。同図に示される構成は上記システムを実現できるコンピュータの一例であり、そのようなコンピュータはこの構成に限定されるものではない。

10

20

30

40

50

## 【 0 0 2 6 】

R O M 3 0 2 は、コンピュータを制御する音声合成プログラムを含む各プログラムを記憶するメモリである。R A M 3 0 3 は、各プログラムの実行時に、R O M 3 0 2 に記憶されているプログラム又はデータを一時的に格納するメモリである。

## 【 0 0 2 7 】

外部記憶装置 3 0 6 は、例えば S S D (ソリッドステートドライブ) 記憶装置またはハードディスク記憶装置であり、入力テキストデータ、入力音声データ、接続音声素片データ、または合成音声データ等の保存に用いられる。また、外部記憶装置 3 0 6 は、図 2 のデータ構成例を有する音声コーパスが格納された音声 D B 1 0 2 を記憶する。

## 【 0 0 2 8 】

C P U 3 0 1 は、各プログラムを、R O M 3 0 2 から R A M 3 0 3 に読み出して実行することにより、当該コンピュータ全体の制御を行う。

## 【 0 0 2 9 】

入力装置 3 0 4 は、ユーザによるキーボードやマウス等による入力操作を検出し、その検出結果を C P U 3 0 1 に通知する。また、入力装置 3 0 4 は、図 1 の入力部 1 0 3 の音声入力装置 1 1 2 の機能を備え、特に図示しないマイクまたはライン入力端子を介して入力音声データを入力し、A / D (アナログ - デジタル) 変換によりデジタルデータに変換した後に、R A M 3 0 3 または外部記憶装置 3 0 6 に記憶させる。さらに、入力装置 3 0 4 は、図 1 の入力部 1 0 3 のテキスト入力装置 1 1 3 の機能を備え、特に図示しないキーボードまたはデバイスインタフェース等を介して入力テキストデータを入力し、R A M 3 0 3 または外部記憶装置 3 0 6 に記憶させる。

## 【 0 0 3 0 】

出力装置 3 0 5 は、C P U 3 0 1 の制御によって送られてくるデータを表示装置や印刷装置に出力する。また、出力装置 3 0 5 は、C P U 3 0 1 が外部記憶装置 3 0 6 または R A M 3 0 3 に出力した合成音声データを、特に図示しないが、D / A 変換器でアナログ合成音声信号に変換した後、増幅器で増幅し、スピーカを介して、合成音声として放音する。

## 【 0 0 3 1 】

可搬記録媒体駆動装置 3 0 7 は、光ディスクや S D R A M、コンパクトフラッシュ等の可搬記録媒体 3 1 0 を収容するもので、外部記憶装置 3 0 6 の補助の役割を有する。

## 【 0 0 3 2 】

通信インターフェース 3 0 8 は、例えば L A N (ローカルエリアネットワーク) 又は W A N (ワイドエリアネットワーク) の通信回線を接続するための装置である。

## 【 0 0 3 3 】

本実施形態による音声合成装置 1 0 0 では、C P U 3 0 1 が、R O M 3 0 2 に記憶された音声合成プログラムを、R A M 3 0 3 をワークメモリとして使用しながら実行することにより、図 1 の音声合成部 1 0 1 内の 1 0 5 ~ 1 1 1 の各ブロックの機能を実現する。そのプログラムは、例えば外部記憶装置 3 0 6 や可搬記録媒体 4 1 0 に記録して配布してもよく、或いはネットワーク接続装置 3 0 8 によりネットワークから取得できるようにしてもよい。

## 【 0 0 3 4 】

図 4 は、図 1 の構成に対応する音声合成装置 1 0 0 の機能を、図 3 のハードウェア構成例を有するコンピュータの C P U 3 0 1 が、ソフトウェアプログラムの処理により実現する場合の、音声合成処理の例を示すフローチャートである。以下随時、図 1、図 2、および図 3 を参照する。

## 【 0 0 3 5 】

C P U 3 0 1 はまず、テキスト入力装置 1 1 3 が入力した入力テキストデータに対して、テキスト解析処理を実行する (ステップ S 4 0 1)。ここでは、C P U 3 0 1 は、入力テキストデータに対して例えば形態素解析処理を実行することにより、入力テキストデータに対応するアクセント付きの音素列を抽出する。この処理は、図 1 のテキスト解

10

20

30

40

50

析モジュール１０５の機能を実現する。

【００３６】

次に、ＣＰＵ３０１は、音声入力装置１１２が入力した入力音声データに対して、韻律解析処理を実行する（ステップＳ４０２）。ここでは、ＣＰＵ３０１は、入力音声データに対して例えばピッチ抽出処理とパワー分析処理を実行する。そして、ＣＰＵ３０１は、ステップＳ４０２のテキスト解析処理により得られたアクセント付きの音素列を参照することにより、音素ごとのピッチの高さ（周波数）、持続時間長、およびパワー（強度）を算出し、それらの情報を目標韻律として出力する。

【００３７】

次に、ＣＰＵ３０１は、素片選択処理を実行する（ステップ４０３）。ここでは、ＣＰ  
Ｕ３０１は、図２に例示されるデータ構成を有する音声コーパスが登録されている音声Ｄ  
Ｂ１０２から、音素および韻律に関して計算されるコストが最良になるように、ステップ  
Ｓ４０１で算出された音素列およびステップＳ４０２で算出された目標韻律に対応する音  
声素片の列を選択する。このとき、ＣＰＵ３０１はまず、音声コーパス中の音素ラベルの  
列（図２（ｂ））をステップＳ４０１で算出された音素列と比較することにより、素片評  
価のコスト条件を満たす素片候補データを音声コーパスからリストアップする。次に、Ｃ  
ＰＵ３０１は、素片候補データにおける音響情報（図２（ｃ））を目標韻律と比較するこ  
とにより、接続評価のコスト条件を満たす最良の素片候補データを、リストアップした素  
片候補データから選択し、最終的に、音声素片の列を選定する。

【００３８】

次に、ＣＰＵ３０１は、波形接続処理を実行する（ステップＳ４０４）。ここでは、Ｃ  
ＰＵ３０１は、ステップＳ４０３での音声素片の選択結果を入力して、対応する音声素片  
の音声データ（図２（ａ））を音声ＤＢ１０２中の音声コーパスから抽出し、それらを接  
続して接続音声素片を出力する。

【００３９】

上述のようにして出力された接続音声素片は、音声ＤＢ１０２が保有している音声コー  
パス内で、入力された音素列と目標韻律に対して、音素に関する素片評価と韻律に関する  
接続評価を合わせたコストが最良になるように選択されたものである。しかし、音声コー  
パスとして巨大なデータベースを保有できないような小規模のシステムにおいては、入力  
音声データから生成された目標韻律と音声コーパス内の限られた規模の音声素片が持つて  
いる韻律とでは、抑揚の付け方などに関して個人に依存する相違点がある。このため、ス  
テップＳ４０４で接続音声素片が出力された段階では、入力音声データで表現されている  
抑揚が接続音声素片に十分に反映されているとはいえない。一方において、単純に目標韻  
律中のピッチやパワーに合うように接続音声素片のピッチやパワーを合わせようすると  
、目標韻律中のピッチやパワーの細かい変動が接続音声素片のピッチやパワーに影響を及  
ぼしてしまい、逆に音質劣化が顕著になってしまう。

【００４０】

そこで、本実施形態では、目標韻律中のピッチやパワーの大局的な変動が発話者の抑揚  
すなわち感情を良く表していると考え、目標韻律からピッチやパワーの緩やかな変動を抽  
出し、その変動データに基づいて接続音声素片のピッチやパワーをシフトさせることによ  
り、目標韻律中に含まれる抑揚情報が良く反映された合成音声を生成する。

【００４１】

そのために、ＣＰＵ３０１は、ステップＳ４０４の波形接続処理の後、ピッチ適応処理  
を実行する（ステップＳ４０５）。図５は、ピッチ適応処理の説明図である。図５（ａ）  
に例示されるように、ＣＰＵ３０１はまず、目標韻律からピッチ周波数の時間変化をピッ  
チ列として抽出する。次に、図５（ｂ）に例示されるように、ＣＰＵ３０１は、ピッチ列  
の各周波数値を適当な粗さで量子化し、量子化されたピッチ列を算出する。これにより、  
目標韻律中の微細なピッチの変動が排除され、ピッチの変化の概形が得られる。さらに、  
図５（ｃ）に例示されるように、ＣＰＵ３０１は、量子化されたピッチ列に対して、時間  
方向の加重移動平均を演算することによって時間方向の平滑化を行い、平滑化されたピッ

10

20

30

40

50

チ列を算出する。具体的には例えば、CPU301は、量子化されたピッチ列で演算中心サンプル位置を先頭から1サンプルずつ移動させながら、その演算中心サンプル位置の両側所定サンプル分について、例えば演算中心サンプル位置から遠ざかるに従って周波数値が一定量ずつ線形に小さくなるようにして、それらの平均値を算出し、その平均値をその演算中心サンプル位置の演算された値として算出する。これにより、図5(a)に例示される細かく変動するピッチ列に対応して、図5(c)に例示されるような自然なピッチ変化を有する平滑化されたピッチ列を得ることができる。CPU301は、このようにして生成した平滑化されたピッチ列の各時間ごとのピッチの値対応するように、ステップS404で出力された接続音声素片の各時間ごとのピッチをシフトさせ、その結果を出力する。

10

#### 【0042】

続いて、CPU301は、ステップS405のピッチ適応処理の後、パワー適応処理を実行する(ステップS406)。なお、ピッチ適応処理とパワー適応処理の実行順番はどちらでもよく、また、どちらか一方のみが実行されてもよい。図6は、パワー適応処理の説明図である。CPU301はまず、図6(a-1)に例示されるように、目標韻律からパワー値の列(以下「パワー列」)を抽出し、同様に、図6(a-2)に例示されるように、接続音声素片(ステップS405のピッチシフトの結果)からパワー列を抽出する。次に、CPU301は、それぞれのパワー列に対して、ピッチ列の場合と同様に時間方向の加重移動平均を演算することにより時間方向の平滑化を行い、図6(b-1)に例示される目標韻律に対応する平滑化されたパワー列と、図6(b-2)に例示される接続音声素片に対応する平滑化されたパワー列を算出する。これにより、それぞれのパワー列において、微細な変動が排除され、パワーの変化の概形が得られる。さらに、CPU301は、目標韻律に対応する平滑化されたパワー列の各時間ごとのサンプル値と、図6(b-2)に例示される接続音声素片に対応する平滑化されたパワー列の各時間ごとのサンプル値との比を算出する。そして、CPU301は、各時間ごとに算出した比の値を、接続音声素片(ステップS405のピッチシフトの結果)の各サンプル値に乗算し、その結果を最終的な合成音声として出力する。

20

#### 【0043】

CPU301は、上述のようにして出力された合成音声データを、例えばRAM303や外部記憶装置306に音声ファイルとして保存するとともに、図1の音声出力装置114を介して合成音声として放音させる。

30

#### 【0044】

図7は、図4のステップS405のピッチ適応処理の詳細例を示すフローチャートである。

#### 【0045】

CPU301はまず、図4のステップS402で生成された目標韻律からピッチ列(以下これを「目標ピッチ列」と記載)を抽出し、この目標ピッチ列と接続音声素片のピッチ列の時間スケールを合わせるタイムストレッチ処理を実行する(ステップS701)。これにより、両者の時間の長さの違いが吸収される。

#### 【0046】

次に、CPU301は、ステップS701でタイムストレッチ処理した目標ピッチ列と接続音声素片のピッチ列のピッチ存在区間を調整する(ステップS702)。具体的には、CPU301は例えば、接続音声素片のピッチ列と目標ピッチ列とを比較し、接続音声素片でピッチの存在しない区間に対応する目標ピッチ列の区間のピッチを削除する。

40

#### 【0047】

次に、CPU301は、ステップS702でピッチ存在区間を調整した後の目標ピッチ列の周波数値を量子化する(図5(b)に対応)(ステップS703)。具体的には、CPU301は例えば、ピッチ周波数値を1オクターブあたりN分割(より具体的には3~10分割等)した単位で、目標ピッチ列を量子化する。

#### 【0048】

50



さらに、CPU301は、ステップS703で量子化した目標ピッチ列を、図5(c)で前述した加重移動平均演算によって平滑化する(ステップS704)。

【0049】

最後に、CPU301は、ステップS704で算出された平滑化された目標ピッチ列を、接続音声素片に適應させる(ステップS705)。具体的には、図5で前述したように、CPU301は、ステップS704で平滑化されたピッチ列の各時間ごとのピッチの値に対応するように、ステップS701で調整された接続音声素片の各時間ごとのピッチをシフトさせ、その結果を出力する。

【0050】

図8は、図4のステップS406のパワー適應処理の詳細例を示すフローチャートである。

10

【0051】

CPU301はまず、図4のステップS402で生成された目標韻律からパワー列(以下これを「目標パワー列」と記載)を抽出し、この目標パワー列と接続音声素片のパワー列の時間スケールを合わせるタイムストレッチ処理を実行する(ステップS801)。なお、図7のステップS701で実行されたタイムストレッチ処理の結果にスケールが合うように調整される。

【0052】

次に、CPU301は、ステップS801でタイムストレッチ処理した目標パワー列と接続音声素片のパワー列のそれぞれを、図6(b-1)および(b-2)で前述した加重移動平均演算によって平滑化する(ステップS802)。

20

【0053】

続いて、CPU301は、ステップS802で算出された目標韻律に対応する平滑化されたパワー列の各時間ごとのサンプル値と、接続音声素片に対応する平滑化されたパワー列の各時間ごとのサンプル値との比を算出する(ステップS803)。

【0054】

最後に、CPU301は、ステップS803で各時間ごとに算出した比の値を、接続音声素片に適應させる(ステップS804)。具体的には、図6で前述したように、CPU301は、ステップS803で各時間ごとに算出した比の値を、接続音声素片の各サンプル値に乗算し、その結果を最終的な合成音声として出力する。

30

【0055】

以上説明した実施形態では、目標韻律中のピッチやパワーの大局的な変動が発話者の抑揚すなわち感情を良く表していると考え、目標韻律からピッチやパワーの緩やかな変動を抽出し、その変動データに基づいて接続音声素片のピッチやパワーをシフトさせることにより、目標韻律中に含まれる抑揚情報が良く反映された合成音声を生成している。しかし本発明では、抑揚情報は目標韻律中のピッチやパワーの大局的な変動に限られるものではない。例えば、抑揚情報として、図4のステップS401で音素列とともに抽出されるアクセント情報を用い、アクセント位置で図4のステップS404の波形接続処理で出力される接続音声素片に対して何らかの加工を行うような適應処理が実行されてもよい。その他、抑揚情報を表現できるようなパラメータを入力音声データから抽出することができれば、そのパラメータによって接続音声素片を加工するような適應処理が実行されてもよい。

40

【0056】

以上のようにして、本実施形態では、波形接続方式の音声合成システムにおいて目標韻律を音声入力によって指定する際に、音声入力による抑揚指定の高い自由度を維持し、かつ音声コーパスの規模を拡大させる必要なく、合成音声の音質を向上させることが可能となる。

【0057】

以上の実施形態に関して、更に以下の付記を開示する。

(付記1)

50

入力されたテキストデータに基づいて、音声素片を選択する選択部と、  
前記選択された音声素片を接続することにより接続音声素片を生成する接続部と、  
前記入力された音声データに含まれる韻律情報から抑揚情報を抽出する抑揚情報抽出部と、

前記接続音声素片に含まれる要素の少なくとも一部を前記抑揚情報に基づいて変更して合成音声を出力する出力部と、

を備える音声合成装置。

(付記 2)

前記出力部は、前記接続音声素片に含まれる要素であるピッチ列を、前記抑揚情報抽出部から抽出された抑揚情報に適應するように変更するピッチ適應部を含む、付記 1 に記載の音声合成装置。

10

(付記 3)

前記抑揚情報抽出部は、前記入力された音声データに含まれる韻律情報としてのピッチ列のピッチを平滑化し、当該ピッチの平滑化されたピッチ列を前記抑揚情報として抽出する、付記 2 に記載の音声合成装置。

(付記 4)

前記抑揚情報抽出部は、前記入力された音声データに含まれる韻律情報としてのピッチ列を構成するピッチを量子化し、当該量子化されたピッチを加重移動平均演算することにより、前記ピッチの平滑化されたピッチ列を生成する、付記 3 に記載の音声合成装置。

(付記 5)

20

前記ピッチ適應部はさらに、前記入された力音声データに含まれる韻律情報としてのピッチ列と前記接続音声素片に含まれるピッチ列との時間スケールを調整するとともに、前記韻律情報としてのピッチ列と前記接続音声素片に含まれるピッチ列とのピッチ存在区間を調整する、付記 2 ないし 4 のいずれかに記載の音声合成装置。

(付記 6)

前記出力部は、前記接続音声素片に含まれる要素であるパワー列を、前記抑揚情報抽出部から抽出された抑揚情報に適應するように変更するパワー適應部を含む、付記 1 に記載の音声合成装置。

(付記 7)

前記抑揚情報抽出部は、前記入された力音声データに含まれる韻律情報としてのパワー列平滑化し、当該平滑化されたパワー列を前記抑揚情報として抽出し、

30

前記パワー適應部は、前記接続音声素片に含まれるパワー列を平滑化し、当該平滑化されたパワー列と前記抑揚情報としての平滑化されたパワー列との比の列を算出し、当該比の列に基づいて前記接続音声素片のパワー列を修正する、付記 6 に記載の音声合成装置。

(付記 8)

前記前記平滑化されたパワー列は、前記パワー列に含まれるパワーそれぞれを加重平均演算することにより取得する、付記 7 に記載の音声合成装置。

(付記 9)

前記パワー適應部はさらに、前記入力された音声データに含まれる韻律情報としてのパワー列及び前記接続音声素片に含まれるパワー列それぞれの時間スケールを調整する、付記 6 ないし 8 のいずれかに記載の音声合成装置。

40

(付記 10)

音声合成装置に用いられる音声合成方法であって、前記音声合成装置が、  
入力されたテキストデータに基づいて、音声素片を選択し、  
前記選択された音声素片を接続することにより接続音声素片を生成し、  
前記入力された音声データに含まれる韻律情報から抑揚情報を抽出し、  
前記接続音声素片に含まれる要素の少なくとも一部を前記抑揚情報に基づいて変更して合成音声を出力する、音声合成方法。

(付記 11)

音声合成装置に用いられる音声合成装置として用いられるコンピュータに、

50

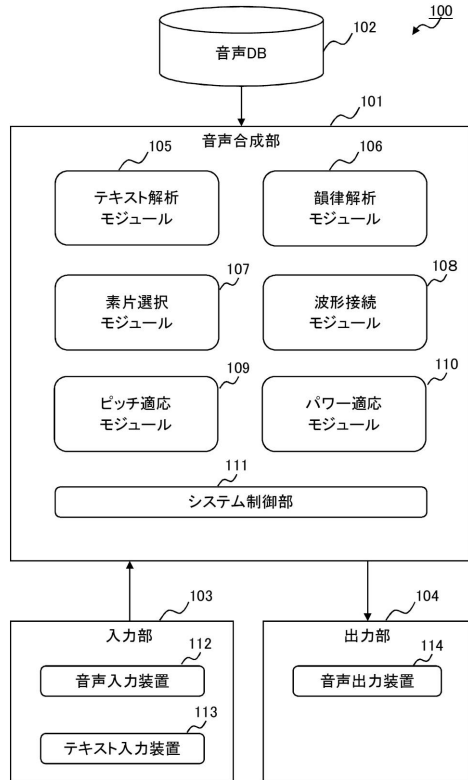
入力されたテキストデータに基づいて、音声素片を選択するステップと、  
前記選択された音声素片を接続することにより接続音声素片を生成するステップと、  
前記入力された音声データに含まれる韻律情報から抑揚情報を抽出するステップと、  
前記接続音声素片に含まれる要素の少なくとも一部を前記抑揚情報に基づいて変更して  
合成音声を出力するステップと、  
を実行させるプログラム。

【符号の説明】

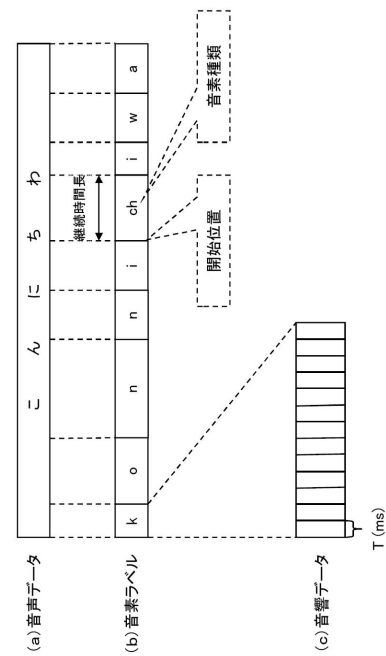
【 0 0 5 8 】

1 0 0	音声合成装置	
1 0 1	音声合成部	10
1 0 2	音声ＤＢ（データベース）	
1 0 3	入力部	
1 0 4	出力部	
1 0 5	テキスト解析モジュール	
1 0 6	韻律解析モジュール	
1 0 7	素片選択モジュール	
1 0 8	波形接続モジュール	
1 0 9	ピッチ適応モジュール	
1 1 0	パワー適応モジュール	
1 1 1	システム制御部	20
1 1 2	音声入力装置	
1 1 3	テキスト入力装置	
1 1 4	音声出力装置	
3 0 1	ＣＰＵ	
3 0 2	ＲＯＭ	
3 0 3	ＲＡＭ	
3 0 4	入力装置	
3 0 5	出力装置	
3 0 6	外部記憶装置	
3 0 7	可搬記録媒体駆動装置	30
3 0 8	通信インタフェース	
3 0 9	バス	
3 1 0	可搬記録媒体	

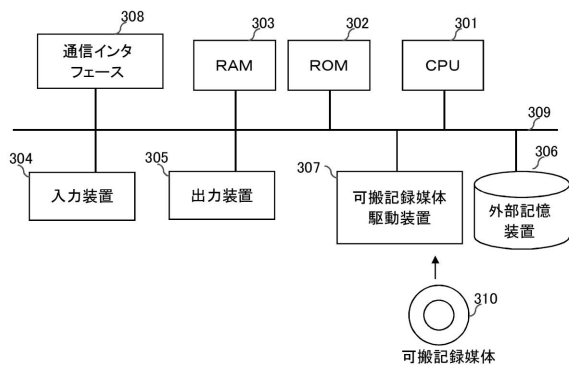
【図 1】



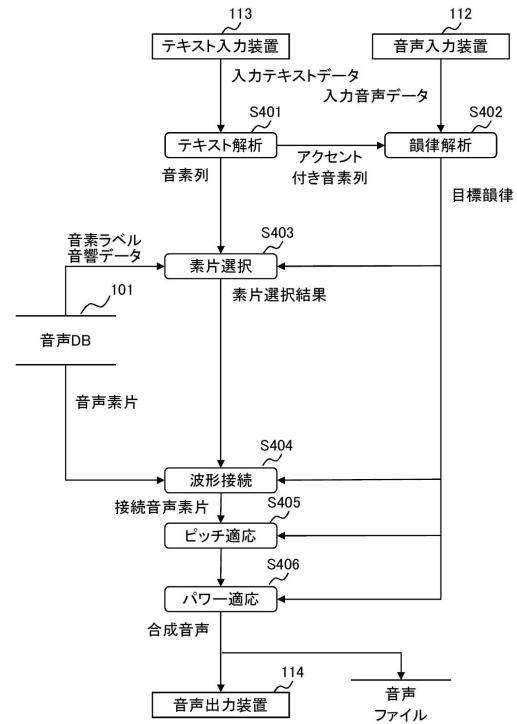
【図 2】



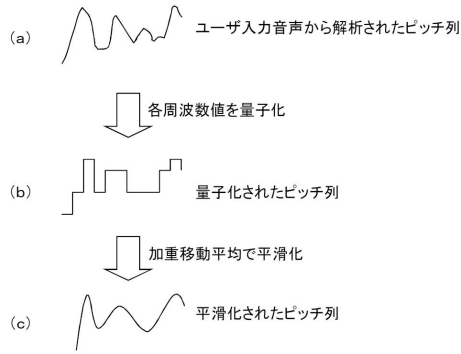
【図 3】



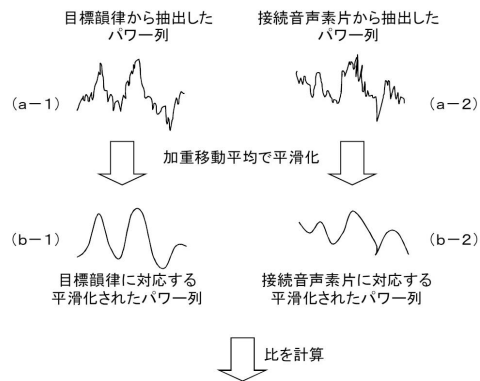
【図 4】



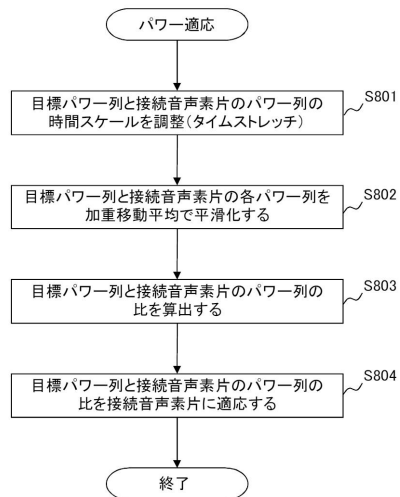
【図 5】



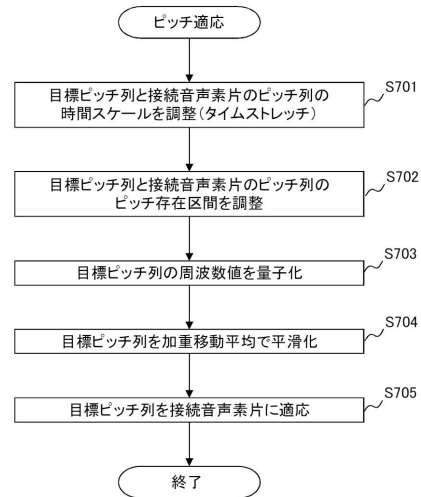
【図 6】



【図 8】



【図 7】



---

フロントページの続き

(56)参考文献 特開2012-220701(JP,A)  
特開2000-010581(JP,A)  
特開2006-309162(JP,A)  
特開平10-153998(JP,A)  
米国特許第05940797(US,A)  
特開2010-009034(JP,A)

(58)調査した分野(Int.Cl., DB名)

G10L 13/00 - 13/10